



LAWRENCE
LIVERMORE
NATIONAL
LABORATORY

Large-Scale Seismic Signal Analysis with Hadoop

T. Addair, D. Dodge, W. Walter, S. Ruppert

October 10, 2013

Computers & Geosciences

Disclaimer

This document was prepared as an account of work sponsored by an agency of the United States government. Neither the United States government nor Lawrence Livermore National Security, LLC, nor any of their employees makes any warranty, expressed or implied, or assumes any legal liability or responsibility for the accuracy, completeness, or usefulness of any information, apparatus, product, or process disclosed, or represents that its use would not infringe privately owned rights. Reference herein to any specific commercial product, process, or service by trade name, trademark, manufacturer, or otherwise does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States government or Lawrence Livermore National Security, LLC. The views and opinions of authors expressed herein do not necessarily state or reflect those of the United States government or Lawrence Livermore National Security, LLC, and shall not be used for advertising or product endorsement purposes.

Large-Scale Seismic Signal Analysis with Hadoop

T. G. Addair₁, D. A. Dodge₁, W. R. Walter₁, S. D. Ruppert₁

₁Lawrence Livermore National Laboratory
7000 East Ave, Livermore, CA 94550, USA

September, 2013

Travis Addair: addair1@llnl.gov
Douglas Dodge: dodge1@llnl.gov
William Walter: walter5@llnl.gov
Stanley Ruppert: ruppert1@llnl.gov

Corresponding author:

Douglas Dodge
Lawrence Livermore National Laboratory
7000 East Avenue
Livermore, CA 94550
MS 046
925-423-4951

38 **Highlights**

- 39 • A global dataset of over 300 million waveforms has been cross correlated.
- 40 • The algorithms have been adapted to run as MapReduce jobs on a Hadoop
- 41 cluster.
- 42 • Increased parallelism was required to make best use of mappers.
- 43 • IO was significantly increased but had little impact on performance.
- 44 • A factor of 19 speedup was achieved relative to initial implementation.
- 45

Abstract

In seismology, waveform cross correlation has been used for years to produce high-precision hypocenter locations and for sensitive detectors. Because correlated seismograms generally are found only at small hypocenter separation distances, correlation detectors have historically been reserved for spotlight purposes. However, many regions have been found to produce large numbers of correlated seismograms, and there is growing interest in building next-generation pipelines that employ correlation as a core part of their operation. In an effort to better understand the distribution and behavior of correlated seismic events, we have cross correlated a global dataset consisting of over 300 million seismograms. This was done using a conventional distributed cluster, and required 42 days. In anticipation of processing much larger datasets, we have re-architected the system to run as a series of MapReduce jobs on a Hadoop cluster. In doing so we achieved a factor of 19 performance increase on a test dataset. We found that fundamental algorithmic transformations were required to achieve the maximum performance increase. Whereas in the original IO-bound implementation, we went to great lengths to minimize IO, in the Hadoop implementation where IO is cheap, we were able to greatly increase the parallelism of our algorithms by performing a tiered series of very fine-grained (highly parallelizable) transformations on the data. Each of these MapReduce jobs required reading and writing large amounts of data. But, because IO is very fast, and because the fine-grained computations could be handled extremely quickly by the mappers, the net was a large performance gain.

69	Keywords
70	• Correlation
71	• Hadoop
72	• MapReduce
73	• Seismology
74	

1. Introduction

It has long been recognized that in collections of seismic data recorded by the same instrument from sources in similar locations there will be many similar seismograms (e.g. Geller and Mueller, 1980, Poupinet et al, 1984). Geller and Mueller (1980) attributed the similarity to the fact that for small magnitude earthquakes, propagation results in an effective low-pass-filtering of the signals to become essentially the Green's functions, so repeated ruptures of the same asperity produce the same seismogram.

Since the initial observations of doublets, researchers have exploited the phenomenon in a variety of different ways. Much work has centered on producing high-precision relocations of clustered seismicity by correlating the waveforms to obtain high-precision relative picks used to relocate the events. For example Fremont and Malone (1987) and Got et al (1994) imaged structures underneath active volcanoes by relocation of multiplets. Rubin et al (1999) imaged seismicity on creeping sections of the Hayward fault. Hauksson and Shearer (2005) relocated 327,000 Southern California earthquakes using waveform cross correlation.

Waveform correlation can also be used as the basis for highly sensitive detectors. This application has been known since at least the 1960's (Anstey, 1966) and has been employed on numerous occasions since. Because the correlation "footprint" of high frequency signals is generally quite small, correlation detectors have been used almost exclusively as "spotlight" detectors aimed at small regions of interest. Harris

and Dodge (2011) used correlation in combination with subspace detectors (e.g. Harris, 2006) in an automated system to track events in an aftershock sequence.

However, interest in using correlation to process events over broader regions has grown. For example Schaff and Richards (2004, 2011) discovered that about 13% of 18,000 earthquakes recorded at regional distances in China were sufficiently well correlated that they could be detected and located using waveform correlation.

To make better use of correlation relationships in seismic data analysis we need to understand what fraction of observed seismicity displays correlation relationships as functions of observed properties (e.g. source-receiver distance, window-length, bandwidth), and correlated source differences (e.g. location, depth, magnitude). The longer-term goal is to understand the physics underlying observed correlation behavior in terms of source similarity and path effects.

However, an impediment to our ability to investigate seismogram correlation is the computational costs of determining these relationships for the tens of thousands of seismic events each year observed at each of the tens of thousands of stations where data is available. Without an extremely efficient way of performing these computations, it is simply too difficult and time consuming to perform the many correlation runs that may be required to achieve a deep understanding of the phenomena under study.

This paper presents a methodology to significantly improve the speed at which massive amounts of seismic event data can be processed to look for correlation behavior using a Hadoop architecture. A result is to explore and expose the correlation behavior of large collections of seismic data. An expected outcome after analysis of these results is completed will be to make correlation a more useful tool in both geophysical research and seismic event cataloging operational systems.

2. Applying correlation to next-generation seismic pipelines

Organizations tasked with monitoring seismicity around the world (e.g. the United States Geological Survey National Earthquake Information Center, the Comprehensive nuclear-Test-Ban Treaty (CTBT) International Data Center, the US National Data Center, the International Seismological Center) and many in specific regions use a processing paradigm that was developed in the 1980's when average computer processing power was a tiny fraction of what is commonly available now. A single pass is made over the waveform data to extract a compact set of parameters such as seismic phase arrival time, amplitude and period. These are input to a phase associator, and the associated phases are fed to a locator, which produces the hypocenter solution. Although numerous refinements have been added over time, the basic procedure is unchanged.

For a variety of purposes, from improving the monitoring of the CTBT, to better characterization of earthquake hazard, to natural resource extraction and reservoir monitoring, it is necessary to detect, locate, and identify seismic events down to very low magnitudes even in the presence of interfering signals. As we near the limits of what can be done using the approach discussed above, pattern-matching-based processing looks increasingly attractive. A correlation between a new signal and a reference one, with a sufficiently high statistic, is at once a detection, location, and identification, assuming those properties are known for the template. A correlation detection and identification can be made with as little as one channel, without needing an associator. Since correlation detectors are much more sensitive than the power detectors used in current systems, it is likely that the detection threshold could be pushed down in all regions for which correlators are available.

Before considering the engineering aspects of a large-scale correlation-based seismic pipeline, it is crucial to understand how effective we can expect it to be. We need to know how much of the Earth's seismicity is correlated and how it is distributed. It may be that an appropriate system could only usefully target specific regions, and therefore its scope should be limited accordingly.

At Lawrence Livermore National Laboratory (LLNL) we operate a research database of seismic events and waveforms for nuclear explosion monitoring and other applications. The LLNL database contains several million events associated with more than 300 million waveforms at thousands of stations (Figure 1). We have

165 correlated the waveforms in this database as a first step towards understanding the
166 global distribution of correlated seismicity and to begin construction of a library of
167 pattern detectors that could be used in a template-based seismic processing
168 pipeline.

169
170 In this exploratory effort we correlate catalog events in a number of specific seismic
171 phase windows (e.g. P, S) and the entire signal length, as well as in a number of
172 frequency bands for each window. For each of the distinct station-channels for
173 which we have waveforms, we find all events whose catalog locations are separated
174 by 50 km or less and with average event-station separation $\leq 90^\circ$. For each pair we
175 apply the processing illustrated in Figure 2. In all ~ 6.8 billion windows are
176 processed.

177 A significant step in processing each window is identification of low SNR windows
178 and artifacts. Low SNR is only a problem because it wastes CPU time, but corrupted
179 and/or bad data (e.g. glitches, dropouts, severe clipping, no recorded signal, etc.)
180 correlate quite well with each other and produce invalid results.

181 To remove artifacts and low SNR seismograms we found it necessary to add a
182 preprocessing step using a decision tree classifier quality control framework to
183 automatically remove the problem signals from the correlation results.

184 Our first full-scale attempt at this process was run on an architecture consisting of 4
185 servers with 44 cores and 613 GB of RAM. All metadata and results were stored in
186 an Oracle database, and the ~ 50 terabytes of waveform data were managed by a 2-

head Hitachi file server. The time required to process the data was about 42 days. Over 370 million correlated waveform segments were found, and these results are still being analyzed.

Because we anticipate repeating this processing with even larger datasets and with variations such as multi-channel correlations, we are interested in reducing the processing time significantly. However, a great amount of effort has already been spent optimizing the workflow for the current architecture, and we think only small improvements are possible. Therefore we decided to move the entire workflow to a MapReduce architecture, which is the subject of the rest of this paper.

3. Hadoop Software Framework

Apache Hadoop is an open source software framework for deploying data-intensive distributed applications. It implements a computational paradigm called MapReduce and the Hadoop Distributed File System (HDFS) derived from Google's MapReduce and Google File System (GFS) respectively (Dean and Ghemawat, 2004; Ghemawat, 2003)

The primary motivation of the MapReduce programming model is to create independent tasks that operate on arbitrary partitions of the input dataset in parallel during the map phase. During the subsequent reduce phase, data that must be rejoined or summarized are shuffled together and processed. When run over a

207 distributed file system such as HDFS with nodes connected to dedicated storage, this
208 framework naturally leverages data locality, whereby tasks running on a particular
209 node process data resident to that node. By moving computation to the data,
210 Hadoop allows exceedingly fast IO and compute performance for highly parallelized
211 algorithms.

212 Hadoop is written in Java, and thus writing MapReduce jobs in Hadoop most
213 commonly involves writing explicit mappers and reducers in Java as well. This
214 deliberate approach to MapReduce programming can encourage the developer to
215 attempt to find a simple mapping of the problem to a single mapper and a single
216 reducer. For many applications, including the waveform correlator, limiting the
217 implementation to one MapReduce job means limiting parallelism where it could be
218 of the greatest benefit. Moreover, the standard Java approach requires significant
219 boilerplate for each job and provides little guidance for chaining jobs into a
220 workflow. Consequently, we explored Pig (Natkovich, 2008) in an effort to find a
221 higher-level tool for applying the MapReduce model.

222 Apache Pig is a platform for creating MapReduce workflows with Hadoop. These
223 workflows are expressed as directed acyclic graphs (DAGs) of tasks that exist at a
224 conceptually higher level than their implementations as series of MapReduce jobs.
225 Pig Latin is the procedural language used for building these workflows, providing
226 syntax similar to the declarative SQL commonly used for relational database
227 systems. In addition to standard SQL operations, Pig can be extended with user-
228 defined functions (UDFs) commonly written in Java. We adopted Pig for our

implementation of the correlator to speed up development time, allow for ad hoc workflow changes, and to embrace the Hadoop community's migration away from MapReduce towards more generalized DAG processing (Mayer, 2013). Specifically, in the event that future versions of Hadoop are optimized to support paradigms other than MapReduce, Pig scripts could take advantage of these advances without recoding, whereas explicit Java MapReduce jobs would need to be rewritten.

3.1 First Proof of Concept

As part of a pilot project facilitated by Livermore Computing (LC), LLNL's institutional high-performance computing group, an effort began to explore porting the waveform correlator to Hadoop. LC's development cluster consisted of 10 Westmere nodes with 12 cores and 96 GB of RAM each. Most importantly for the IO bound waveform correlator, each node also featured dedicated storage to promote data locality. To compare the performance of the existing solution to Hadoop, a 1 terabyte (TB) test dataset was derived from a seismically dense region of the western United States and copied into HDFS.

Taking full advantage of the increased computational power and IO throughput of a Hadoop cluster requires significant parallelism of the application-level algorithms. Because every STA-CHAN in the original waveform correlator implementation was an independent task, and there were over 10,000 such tasks to perform, it would seem that the work was more than distributable enough to keep 10 hard-drives and 120 compute cores busy. In practice, however, there was significant imbalance in

the effort performed by each task, as illustrated by Figure 3, which shows the original implementation's time to completion for each task on the 1 TB test dataset. As the histogram shows, most tasks in the original implementation executed in a matter of seconds. For these tasks, IO is the natural bottleneck as they do very little beyond reading a waveform, rejecting it, and requesting the next candidate. However, a handful of STA-CHANs containing dense clusters of high-quality data required several hours to fully process, imposing clear computational constraints on the performance. Any single job MapReduce implementation would be similarly constrained by these outliers, given that each mapper and reducer is run within its own single thread of execution. In order to increase the parallelism of the Hadoop correlator, and consequently improve performance, the problem needed to be broken up into finer-grained sub-problems that could be more evenly distributed across the cluster topology.

3.2 Overview of Pig Workflow

The IO limitation of the original waveform correlator hardware lent itself to an architecture that attempted to minimize IO, going so far as to ensure that each waveform was read no more than once. In contrast, the Hadoop implementation takes advantage of the improved IO throughput by trading reads and writes to gain increased parallelism. As Figure 4 shows, the first proof of concept Hadoop implementation added many additional reads, writes, and transformations to the data that didn't exist in the original implementation. Despite this increase in design complexity, IO, and computation, a dramatic increase in parallelism across each of

the steps allowed us to take full advantage of the cluster hardware. This increased problem granularity meant reduced work done in the bottlenecks, and spread the most computationally intensive procedures over many more threads of execution. A discussion of the major processing steps implemented in this pipeline follows.

3.2.1 Join Metadata

For the purpose of minimizing disk usage and promoting consistency, LLNL maintains a highly normalized relational database management system (RDBMS) for storing seismic STA-CHAN and waveform metadata. This solution has worked well for our more typical interactive, single-workstation use cases, but creates an immediate bottleneck for highly parallel applications. Specifically, any application spread across tens of nodes and hundreds of cores can quickly saturate both the database's CPU and the bandwidth of the network interconnect. Going against conventional RDBMS^[D1] wisdom, it makes sense in an environment such as Hadoop to denormalize the dataset up front. By doing so, the need to request data randomly is eliminated, and the dataset can easily be partitioned and sent to independent tasks for processing. In contrast, a normalized solution would require each task to make requests for missing chunks of metadata, imposing a network limitation on the task.

Pig provides a native "join" function for taking two structured datasets and denormalizing them by some join predicate. Using Apache Sqoop, a command-line tool for transferring data between relational databases and Hadoop, we ingested all necessary Oracle tables into HDFS, and then used Pig to join the datasets into one

294 denormalized `metadata[D2]` table. Fully denormalized, our metadata took just over
295 50 GB of usable space on disk.

296 We also joined the metadata to the binary waveform signal data directly by
297 attaching a STA-CHAN-EVENT key onto the waveform data itself. By defining a
298 simple schema consisting only of two fields (the unique key and the binary
299 waveform data), we were able to perform a similar join in Pig to create a fully
300 denormalized table containing all the necessary data for a given waveform in a
301 single row along with the trace itself.

302 **3.2.2 Get Candidates**

303 Having the dataset fully denormalized into rows of waveforms along with all
304 relevant STA-CHAN and event metadata, we next approached the issue of cutting
305 and filtering the waveforms into valid candidates for correlation. This was in
306 contrast to the original processing pipeline that took raw waveforms, checked for
307 nearest neighbors, then cut and filtered the set of neighbors together. In a scenario
308 where IO is the primary performance consideration, it makes sense to avoid cutting
309 and filtering as a preprocessing step, as it necessarily requires reading all raw
310 waveforms once, then writing and reading back in all cut and filtered waveforms in
311 a second pass.

312 But, when IO is cheap, additional passes over the data to massage them into a more
313 workable state can offer significant parallelization benefits. Cutting and filtering a
314 waveform is an independent operation, and can thus be parallelized more finely
315 than the complete set of operations on a STA-CHAN. For the 1 TB test dataset

316 consisting of nearly 10 million traces on about 10 thousand STA-CHANs, that
317 represents a factor of 1,000 more tasks to be spread across the cluster. Such smaller
318 tasks are necessarily less variable in the amount of work to be done, and thus
319 distribute the computational effort more evenly.

320 Another benefit of performing the candidate generation task early in the processing
321 is that it creates an opportunity to filter the input dataset down to something
322 smaller and more manageable. This may seem counter-intuitive considering that S
323 waveforms cut into W windows and filtered into B bands yields as many as SWB
324 candidates. In practice, however, the majority of these preliminary candidates are
325 discarded for having low signal-to-noise ratios or otherwise failing one of the
326 decision tree classifiers used to assess the quality of the data. For the 1 TB test
327 dataset, only about $\frac{1}{2} S$, or 4.5 million candidates were output from the candidate
328 generation task.

329 **3.2.3 Nearest Neighbors**

330 Left with only the cut and filtered waveforms that passed the initial quality control
331 tests, the nearest neighbor calculation, too, can be parallelized beyond the STA-
332 CHAN granularity with fewer events per task. This is because we only wish to
333 correlate waveforms that are cut into the same window and filtered into the same
334 band, and many of the events for a particular window and band were discarded in
335 the previous step. Thus we can have a separate task for every STA-CHAN-WINDOW-
336 BAND with only a fraction of the events to be processed per task.

337 An important implementation detail of the nearest neighbor calculation is that every
338 task requires all of its event spatial data (latitude/longitude coordinates) to be in
339 memory at once. At this point in the processing, this event metadata is currently in
340 a denormalized table with all the binary waveform data, which will certainly not all
341 fit in memory at once for a given STA-CHAN-WINDOW-BAND. Fortunately, Pig
342 provides a mechanism to specify a subset of the fields (or columns, in relational
343 database terms) to be sent to a given function. To reconcile the fact that the spatial
344 event data are small enough to fit in memory, but the waveform data are not, we
345 calculate the nearest neighbors only once.

346 The output from the nearest neighbor calculation will provide two additional pieces
347 of metadata for each waveform: an “island” identifier and a step number. An island
348 of events is a connected component in graph terminology. Given our criterion that a
349 neighboring event is within 50 km, we can say that no two events in separate islands
350 are[TA3] within 50 km of each other. This is useful for the purpose of defining finer-
351 grained sub-problems, as individual STA-CHAN-WINDOW-BAND-ISLANDs can now
352 be correlated entirely independently, in parallel. However, even islands of events
353 can be too large to fit entirely in memory, and so we need one additional piece of
354 information to constrain the problem.

355 As mentioned previously, we output a step number for each waveform, representing
356 the order in which we traversed the waveform’s event within its island in the
357 nearest neighbor calculation. We traverse the events within a task by first drawing
358 an event from the set at random. We then use an R-Tree (Guttman, 1984) to fetch

all neighboring events within 50 km with logarithmic asymptotic time complexity, and add the neighboring events to a queue. The events on the queue are the next to be processed in a similar manner, adding their previously unseen neighbors to the queue until no more neighbors exist, at which point we have discovered an island. Concretely, we say that for an island of N events, the waveform with step 1 corresponds to the first event we draw at random to seed the island, and the waveform with step N corresponds to the last event drawn from the queue that has no previously unseen neighbors. By processing events in a queue, instead of simply drawing them at random, we perform a breadth-first traversal of the island's corresponding graphical representation. This enables us to process neighboring events together, keeping them in memory while we process their neighbors, then discarding them as we migrate towards sections of the island farther away. Outputting the step value enables us to recreate this migratory processing without needing to keep all the spatial event information in memory during future processing tasks.

With the island and step metadata generated, this information is then rejoined to the complete waveform dataset, and grouped together by STA-CHAN-WINDOW-BAND-ISLAND using standard Pig functions.

3.2.4 Correlate

Every STA-CHAN-WINDOW-BAND-ISLAND can be processed independently. First, all events within an island are ordered by ascending step number. This ordered bag of waveform data is then fed to a user-defined function: CORRELATE. This is a

381 special kind of user-defined function called an “accumulator” that processes
382 elements within the ordered bag individually as opposed to loading them into
383 memory at once. In this way, the CORRELATE function is stateful, invoked with one
384 waveform at a time, and outputting to a separate bag of correlations as neighboring
385 events are fed in. These correlations form the final output of the pipeline once
386 flattened into a conventional comma-delimited text file with appropriate metadata
387 attached.

388 Before we accumulate our first waveform, we define an event queue we will add to
389 in the same order as we accumulate events, and we set a variable “current” to define
390 the waveform against which all incoming waveforms will be correlated until they
391 are no longer neighboring events, at which point we can safely discard the current
392 event. This works because of the way we have ordered our events by step number,
393 such that if an event follows another, either it is a neighbor or there are no more
394 neighbors to process. A formal proof of this claim, with accompanying pseudo-code
395 of the CORRELATE algorithm, can be found in the Appendix.

396 Once there are no more neighbors to process, we remove the next event from the
397 queue and declare it to be the current event. Necessarily, the incoming event must
398 be a neighbor of the newly declared current event, or there are no unseen events for
399 the current event. This is shown by Lemma 1 (Appendix). Once there are no more
400 incoming events to accumulate, we drain the queue of its remaining events:
401 correlating each newly removed event against all remaining events on the queue. In
402 practice, we again speed up the neighbor calculation with an R-Tree.

3.3 Performance Improvements and Bottlenecks

Implementing the Map Reduce workflow described above and deploying it on the Hadoop test cluster yielded a significant performance improvement over the original implementation of the correlator. However, certain bottlenecks in the processing pipeline remained. As Figure 5 illustrates, processing time for 99% of the 1 TB test dataset went from nearly 28 hours on existing hardware to 45 minutes on Hadoop. However, we observe that for the last 1% of the 1 TB test dataset the performance gap does not close as much between the two implementations: 30 hours on existing hardware to over 6 hours on Hadoop.

Our initial suspicion was that the loose connectivity constraint imposed on the islands was leading to sprawling, yet sparse, islands of events that could be broken up into smaller, and minimally overlapping components. However, experimentation with algorithms including minimum graph cuts (Hao and Orlin, 1994) and the density-clustering DBSCAN (Ester et al., 1996) showed that our dataset contained a handful of large, highly dense islands of quality events. This made the prospect of further refining the granularity of our tasks difficult, since a fully connected graph of events requires that every event be correlated against every other event. In other words, in places where seismicity concentrated there can be a very large number of events within 50 km of each other that require all possible pairs of correlations to be generated.

In graph theory terms we can think of an island of neighboring events as a connected component of V vertices (events) and E edges (neighboring relations). In

the CORRELATE process described above, we take time-domain waveform segments, and align them with their neighboring event segments before converting them into the frequency domain via a Fast Fourier Transform (FFT). Once the event and all its neighbors are in the frequency domain, we correlate them to produce a correlation coefficient between 0 and 1. Thus we must perform $O(E)$ FFTs per STA-CHAN-WINDOW-BAND-ISLAND task. For the extremely dense outlier islands mentioned above, $E \approx V^2$, which implies that $O(E) = O(V^2)$. For our test dataset, the largest complete island consisted of about 10,000 events, and thus roughly 100,000,000 FFTs confined to a single thread of execution. Naturally, this led us to consider ways to do better.

3.4 Refined Proof of Concept

The goal we had in mind for improving the performance of CORRELATE was to reduce the $O(V^2)$ FFTs down to a more manageable $O(V)$. An ideal solution would be to simply compute all FFTs once per event (vertex) and use the transformed, frequency domain segments in the CORRELATE step. This solution would have the added benefit of being embarrassingly parallel, capable of being performed as a pre-processing step with one task per waveform segment. However, calculating the FFTs once per event requires additional information about all the events to be correlated that is missing a priori. Before we can transform the segmented data into the frequency domain we need to know the minimum pre- and post-picked arrival times for which all segments to be correlated together have data. Once this information is obtained for a given STA-CHAN-WINDOW-BAND-ISLAND, all the

447 segments can be trimmed down to the same length in seconds. Without trimming
448 the segments to the intersection of their lengths in this manner, we run the risk of
449 producing separate power-of-two length output segments from the FFTs. For the
450 frequency-domain multiplication to work, the basis must of course be the same, and
451 so the minimum pre- and post-pick times must be known prior to performing the
452 FFT.

453

454 Under the constraint that a given waveform could only be read from disk once,
455 calculating island-wide pre- and post-pick times would be impossible without
456 exhausting heap memory, so in our first implementation we opted instead for
457 performing more FFTs in exchange for less IO. However, that implementation
458 revealed that reading and writing the entire dataset to HDFS could be performed on
459 the order of seconds and minutes. Consequently, our revised solution was simple:
460 add two more passes over the data to the pipeline.

461

462 The first new task would accumulate an entire STA-CHAN-WINDOW-BAND-ISLAND
463 and calculate the pre- and post-pick times (GetBounds). Those two pieces of
464 additional metadata would then be appended onto the segment-level metadata, and
465 all FFTs would be calculated (`FourierTransform[D4]`). Each task would take a single
466 waveform and produce exactly one transformed segment, allowing this step to be
467 massively parallelized. This revised workflow is shown in Figure 6. Figure 7 shows
468 that adding the two additional passes over the data, in spite of the increased IO cost,

yielded a factor of 10 performance improvement in the CORRELATE routine. The added cost incurred by the highly dense outlier islands was greatly diminished. Both of the additional processing steps contributed negligibly to the overall runtime of the system (less than 15 minutes combined). Most importantly, the CORRELATE runtime was reduced to under an hour of processing. In aggregate, the refined Hadoop implementation yielded a factor of 19 improvement over the original waveform correlator, going from 48 hours on the 1 TB test dataset to under 3 hours in total[D5].

These performance gains include the time to read and write the data from and to HDFS, and were obtained in spite of the dramatic increase in total IO over the original implementation shown in Figure 8.

4.0 Discussion

The Hadoop model of distributing the data with the computations presents a paradigm shift for the data-intensive scientific computing community. It demonstrates the need to change algorithmic priorities to fully take advantage of these powerful systems. Instead of asking ourselves how we can decrease the IO burden, as we did in our original implementation of the waveform correlator, we now find ourselves asking how we can increase the parallelism of our algorithms. Whereas before we had lots of CPU which could not be fully utilized, now we have blazingly fast IO and imbalanced CPU load. The process of finding new ways to break apart one's algorithms into finer-grained sub-problems is at the heart of the Hadoop philosophy: scale out, not up.

Based on our 1TB test data set and the factor of 19 performance increase found by moving to the HADOOP architecture, we expect we will be able to re-correlate our entire ~50 TB, ~300 million waveform database in about 2 days instead of the original 42 days. This will dramatically improve our ability to conduct research on massive seismic datasets, and we intend to describe those results in future papers. The lessons of this study, making use of HADOOP to increase parallelism instead of reducing IO, apply to many massive datasets of time series data, which are common in geophysics and other fields.

5.0 Acknowledgements

This work was performed under the auspices of the U.S. Department of Energy by Lawrence Livermore National Laboratory under contract DE-AC52-07NA27344. Lawrence Livermore National Security, LLC

References

- Anstey, N.A., 1966. Correlation Techniques—A Review, *Can. J. Expl. Geophys.*, 2, 55–82.
- Dean, J., and S. Ghemawat. MapReduce: Simplified Data Processing on Large Clusters. In OSDI'04, 6th Symposium on Operating Systems Design and Implementation, Sponsored by USENIX, in cooperation with ACM SIGOPS, pages 137–150, 2004.
- Ester, M., H. Kriegel, J. Sander, X. Xu (1996). "A density-based algorithm for discovering clusters in large spatial databases with noise". In Evangelos Simoudis, Jiawei Han, Usama M. Fayyad. *Proceedings of the Second International Conference on Knowledge Discovery and Data Mining (KDD-96)*. AAAI Press. pp. 226–231.
- Fremont, M.J., and S. D. Malone (1987). High precision relative locations of earthquakes at Mount St. Helens, Washington, *J. Geophys. Res.* 92, 10,233–10,236.
- Geller, R. J., and C. S. Mueller, Four similar earthquakes in central California, *Geophys. Res. Lett.*, 7, 821-824, 1980
- Ghemawat, S., H. Gobioff, S. Leung, (2003), "The Google File System", *19th Symposium on Operating Systems Principles* (conference), Lake George, NY: The Association for Computing Machinery.
- Got, J. L., J. Frechet, and F. W. Klein (1994). Deep fault plane geometry inferred from multiplet relative relocation beneath the south flank of Kilauea, *J. Geophys. Res.* 99, 15,375–15,386.
- Guttman, A. (1984). "R-Trees: A Dynamic Index Structure for Spatial Searching". *Proceedings of the 1984 ACM SIGMOD international conference on Management of data - SIGMOD '84*. p. 47. doi:10.1145/602259.602266.
- Harris, D. (2006), Subspace detectors: Theory, Lawrence Livermore Natl. Lab. Rep. UCRL-TR-222758, 46 pp., Lawrence Livermore Natl. Lab., Livermore, Calif.
- Harris, D. and D. Dodge (2011). An autonomous system for grouping events in a developing aftershock sequence, *Bull. Seism. Soc. Am.* 101, 763-774, doi:10.1785/0120100103.
- Hauksson, E., and P. Shearer (2005). Southern California hypocenter relocation with waveform cross-correlation, part 1: Results using the double-difference method, *Bull. Seism. Soc. Am.* 95, 896–903.

Hao, J. X., J., B., Orlin, (1994), "A Faster Algorithm for Finding the Minimum Cut in a Directed Graph". *Journal of Algorithms* 17 (3): 424. doi:10.1006/jagm.1994.1043

Mayer, C., (2013), Hortonworks announce Stinger to solve Hadoop's real-time headache, <http://jaxenter.com/hortonworks-announce-stinger-to-solve-hadoop-s-real-time-headache-46261.html>.

Natkovich, O. (2008), Pig – The Road to an Efficient High-level language for Hadoop, Hadoop Blog. <http://developer.yahoo.com/blogs/hadoop/pig-road-efficient-high-level-language-hadoop-413.html>.

Poupinet, G., W. L. Ellsworth, and J. Frechet (1984). Monitoring velocity variations in the crust using earthquake doublets: an application to the Calaveras fault, California, *J. Geophys. Res.* 89, 5719–5731.

Rubin, A. M., D. Gillard, and J.-L. Got (1999). Streaks of microearthquakes along creeping faults, *Nature* 400, 635–641.

Schaff, D. P., and P. G. Richards (2004). Repeating seismic events in China, *Science* 303, 1176–1178.

Schaff, D. P., and P. G. Richards (2011). On finding and using repeating seismic events in and near China, *J. Geophys. Res.* 116, doi:10.1029/2010JB007895

Figure Captions

Figure 1 shows the waveform density (number of waveforms in database per cell divided by the total number of waveforms). Color is proportional to $\log(\text{density})$ with black lowest and white highest. Note that although the data set has global coverage, the density is highest in the Middle East, Eurasia, and Western North America.

Figure 2 is a schematic illustration of the processing applied to a single channel for a pair of events observed by a single station. For each of B bands the seismograms are filtered and cut into W phase windows. For each window pair, the cross correlation function is computed and the max and its associated shift are recorded in the database.

Figure 3 shows the time to completion for each STA-CHAN task on the original architecture for the test dataset. Outlier tasks take several orders of magnitude longer to complete than average.

Figure 4 shows the first proof of concept implementation of the waveform correlator as a Pig workflow consisting of many finely-grained passes over the data.

Figure 5 is a comparison of processing times for the test dataset on the original architecture (left) and the first Hadoop implementation (right). Times are in seconds.

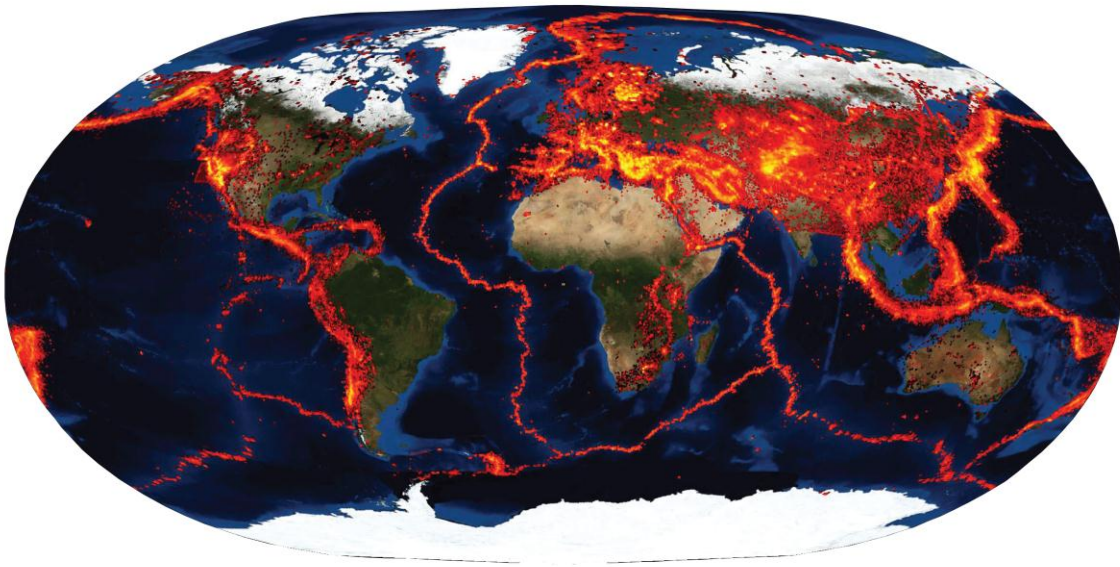
Figure 6 shows the revised Hadoop processing flow with the added “Get Bounds” and “Fourier Transform” processing steps.

Figure 7 shows the comparison of processing times for the test dataset on the original architecture (left) and the final Hadoop implementation (right). Times are in seconds.

Figure 8 is a comparison of read/write times for the original implementation and the final Hadoop implementation. Times are in seconds.

603 **Figures**

LLNL Research Database of Seismic Waveforms



604

605

Figure 1

606

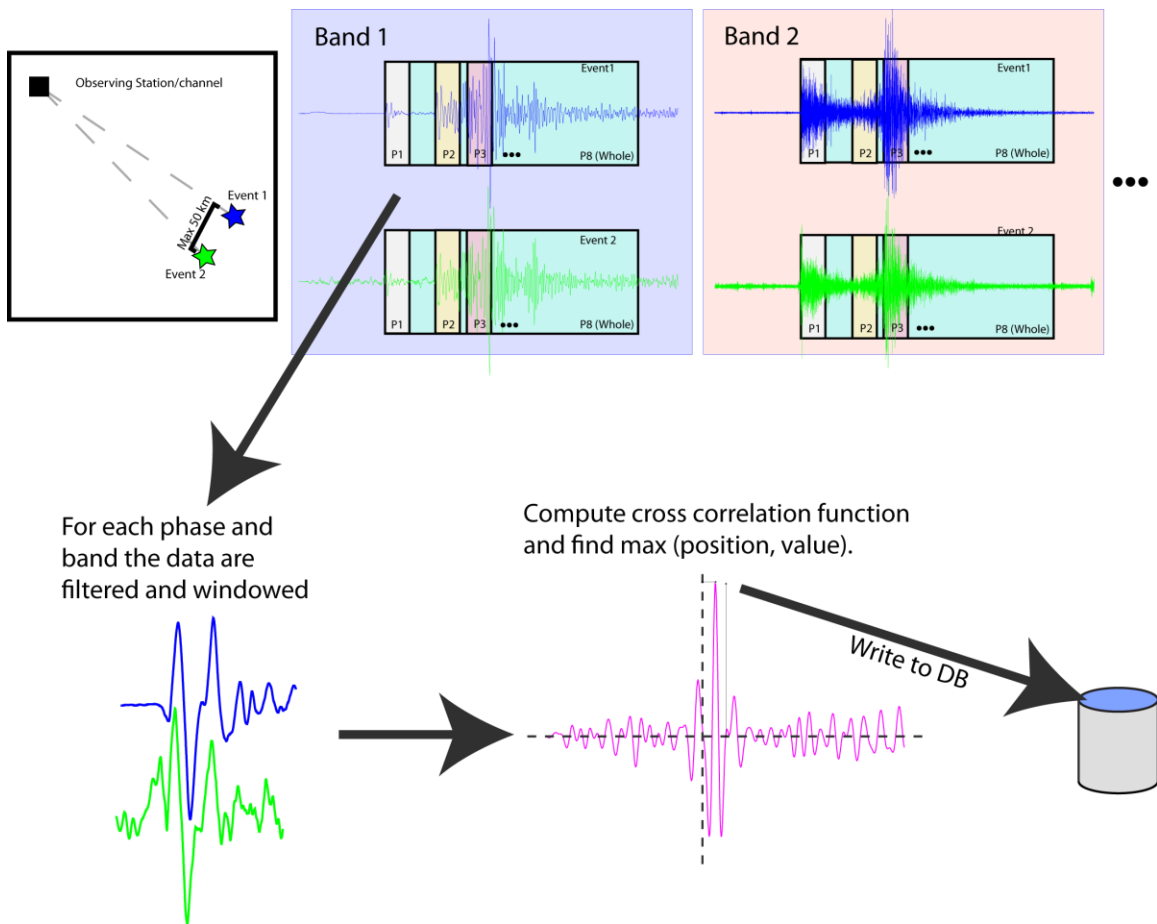
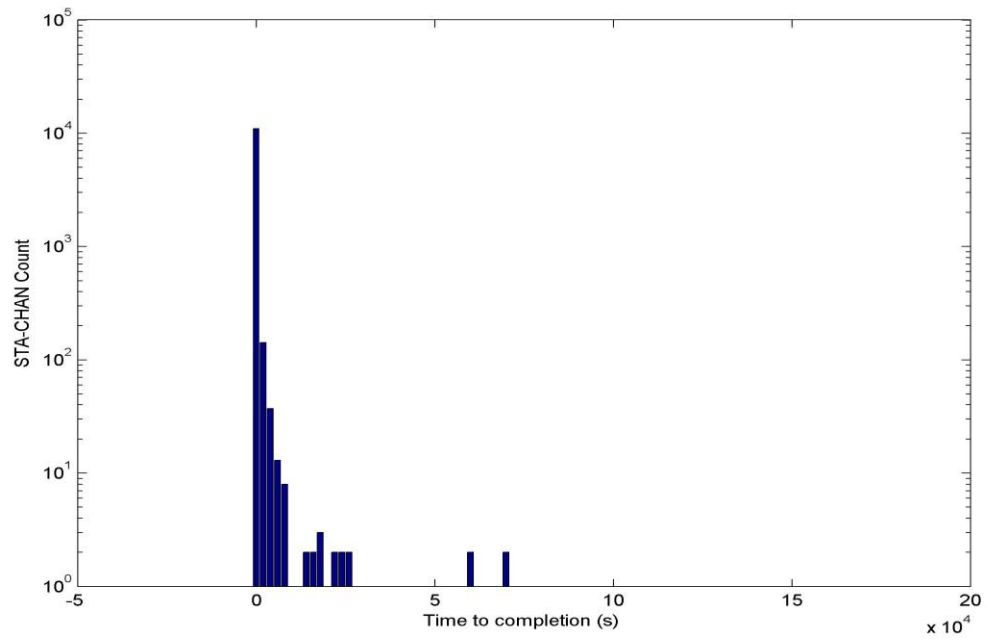


Figure 2

610



611

612

Figure 3

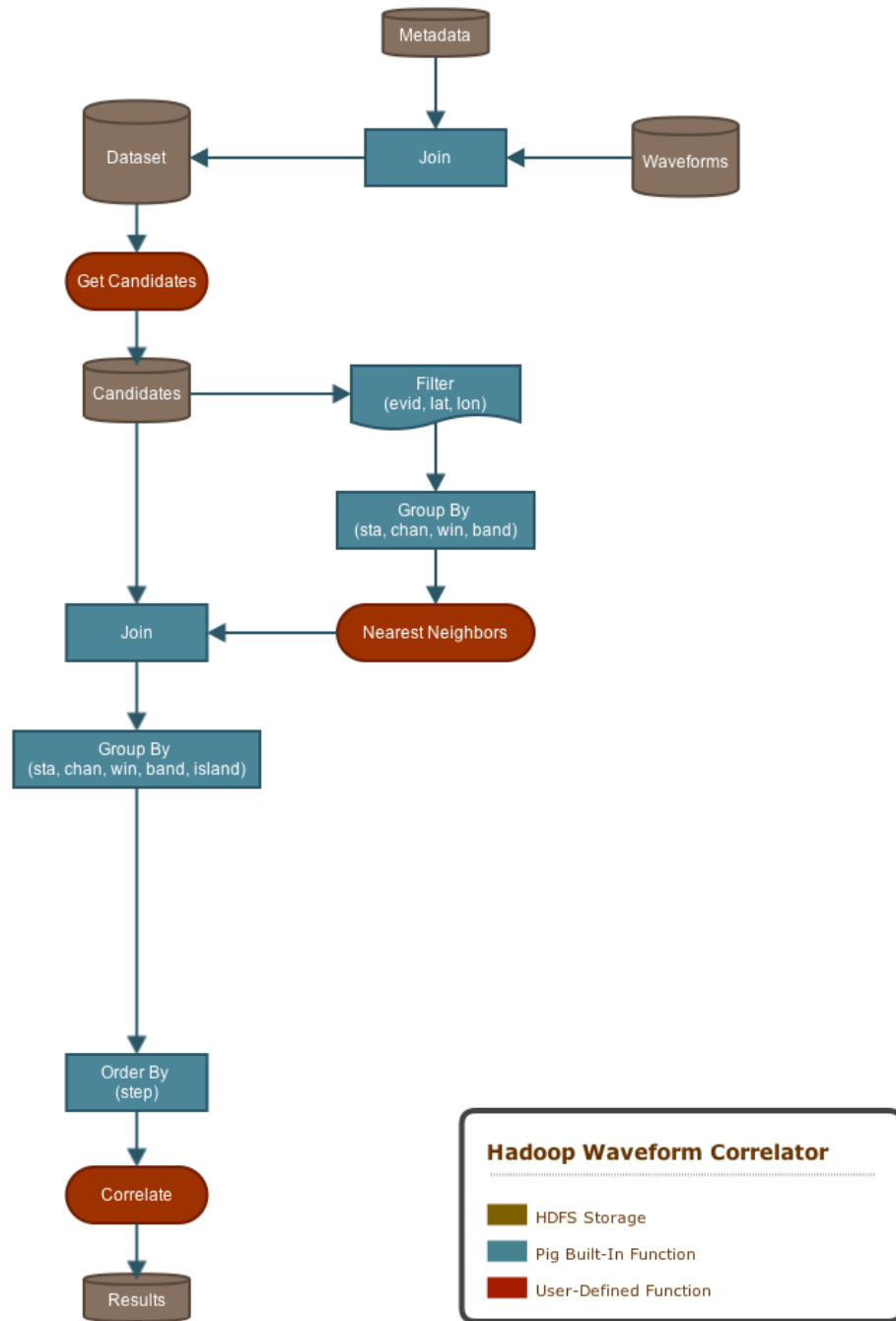


Figure 4

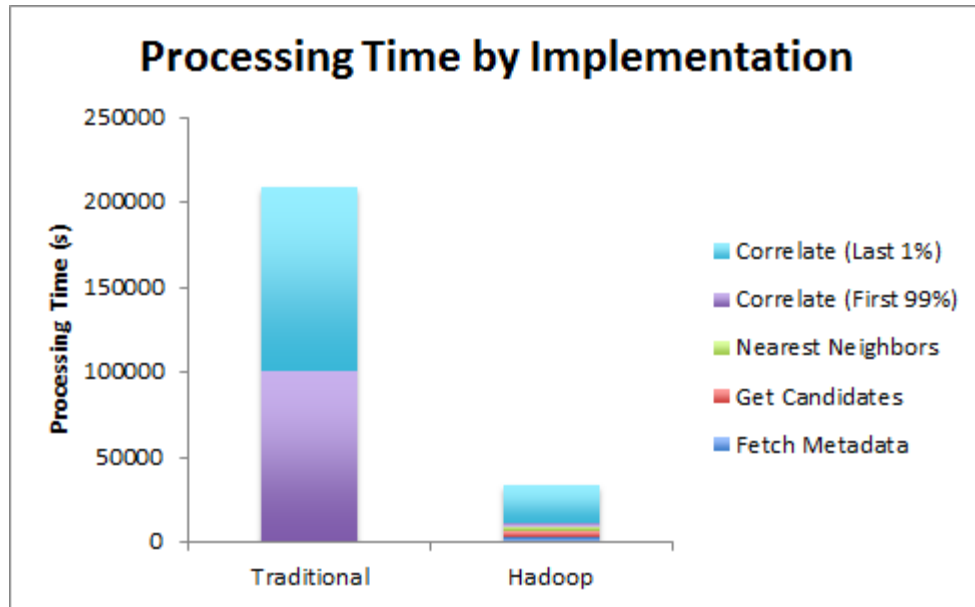


Figure 5

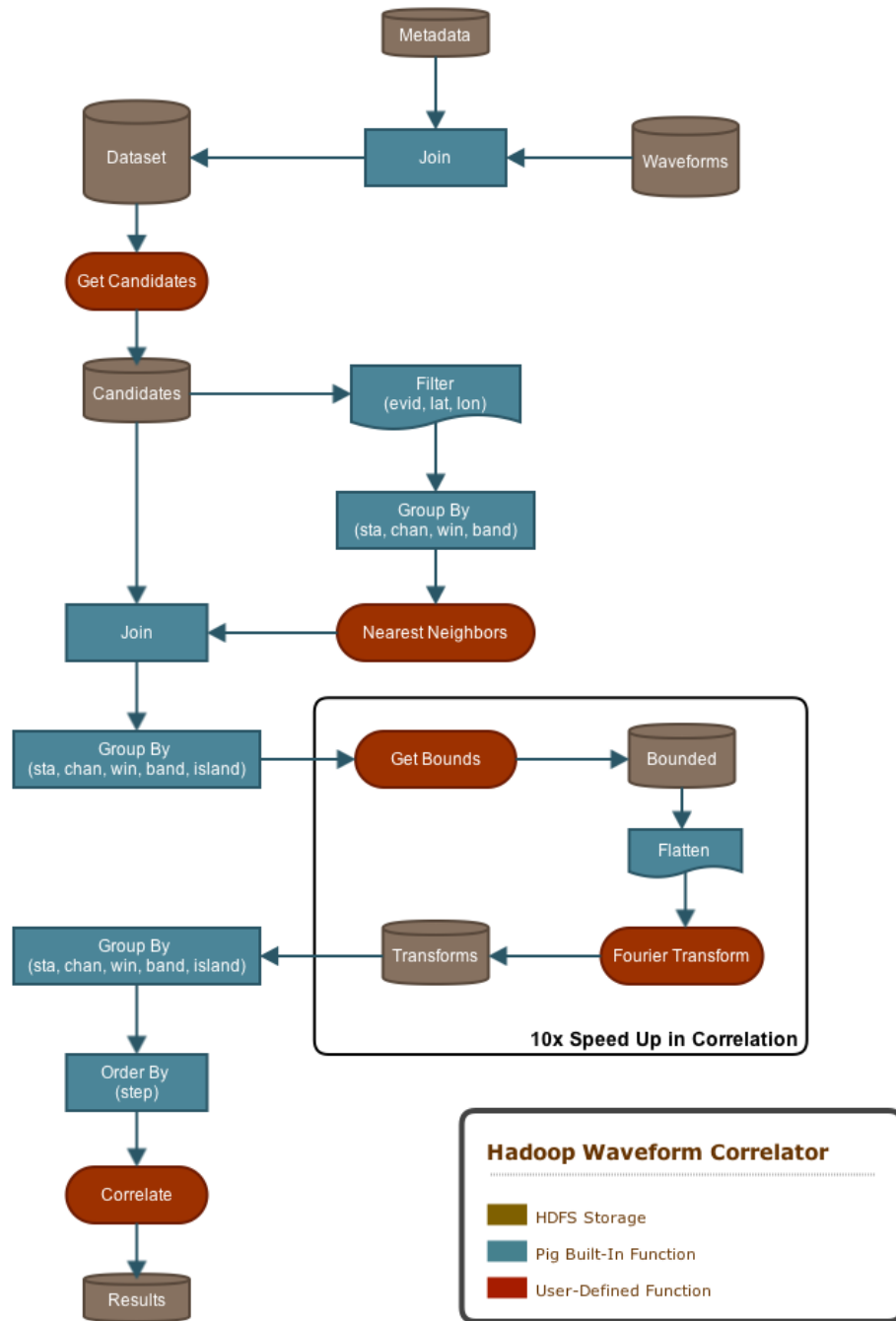


Figure 6

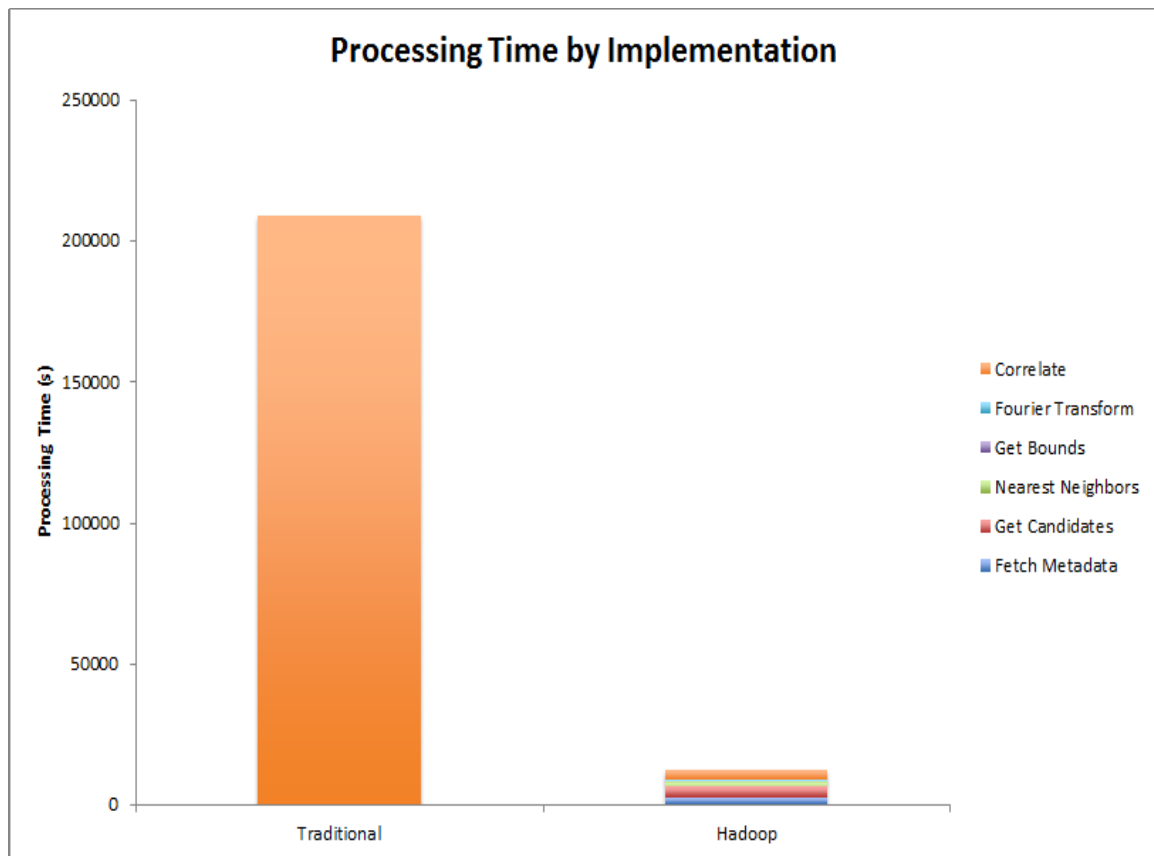


Figure 7

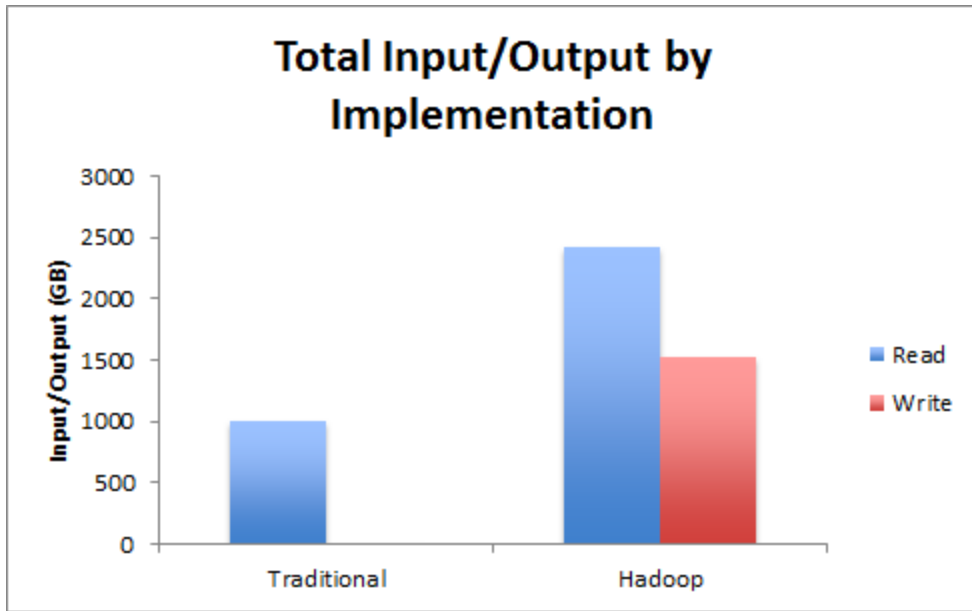


Figure 8

629 **Appendix**

630 The CORRELATE function, in pseudo code, works as follows:

```
631     current = null
632     queue = new Queue()
633
634     CORRELATE(incoming):
635         if current == null:
636             current = incoming
637         else:
638             queue.add(incoming)
639             if neighbors(current, incoming):
640                 correlate(current, incoming)
641             else:
642                 while not neighbors(current, incoming):
643                     current = queue.dequeue()
644                     for w in queue:
645                         if neighbors(current, w):
646                             correlate(current, w)
```

647

648 Lemma 1: In the CORRELATE calculation, there is at all times a **current event** being
649 correlated against all **incoming events** being removed from the **queue**. If an
650 incoming event is not the neighbor of the current event, then all of the current
651 event's neighbors have been correlated against the current event.

652 Proof: In the NEAREST NEIGHBORS calculation, there is at all times a **current event**
653 whose step number corresponds to the order it was removed from the **queue** in the
654 breadth-first search. All of the current event's unseen neighbors are added to the
655 queue together and assigned a **step** number greater than that of the current event.
656 All previously unseen neighbors will be processed in a sequence together called the
657 **unseen sequence**.

658 By induction on the current event in the CORRELATE calculation, we will show that
659 CORRELATE has correlated all of the current event's neighbors (the **current**
660 **candidates**) by the time an incoming event is accumulated that is not a current
661 candidate.

662 In the base case, at step 1, we accumulate the first incoming event from the bag of
663 island events ordered by step number and promote it to be the current event. There
664 are two possibilities: either the current event has an unseen sequence of current
665 candidates to correlate, or it does not. If there is no remaining unseen sequence to
666 process, then our claim is correct by definition. Thus we assume that the current
667 event has an unseen sequence that needs to be processed. If the next incoming
668 event is not a current candidate, then it must be part of another event's unseen
669 sequence by definition of the connectivity of an island. By definition of the current
670 event, the incoming event must be part of the unseen sequence of an event with a
671 higher step number than the current event. But the current event's unseen
672 sequence must necessarily come before the unseen sequence of any subsequent
673 events by definition of the step number, and so this is a contradiction. Thus there
674 are no more neighbors to correlate against the current event.

675 Suppose the claim holds for current events up to step $k-1$. At step k , we remove the
676 k^{th} event accumulated from the queue and promote it to be the current event. We
677 then proceed to correlate it against all the elements on the queue. If the current
678 event had any neighbors with a lower step number than itself, they were correlated
679 earlier in the process by our assumption. Thus the only current candidates

680 remaining must come from incoming events. Without loss of generality, we can
681 apply the same reasoning to incoming events as applied in the base case to
682 demonstrate that all remaining neighbors must be part of the next incoming
683 sequence, or there are none left to correlate against the current event, and so the
684 CORRELATE algorithm does not miss any potential correlations.